| | | |
|---|---|---|
| **Unit 1:** Exploring One-Variable Data | *E* | **15–23%** |
| **Unit 2:** Exploring Two-Variable Data | | **5–7%** ✗ |
| **Unit 3:** Collecting Data | *D* | **12–15%** |
| **Unit 4:** Probability, Random Variables, and Probability Distributions | *C* | **10–20%** |
| **Unit 5:** Sampling Distributions | | **7–12%** ✗ |
| **Unit 6:** Inference for Categorical Data: Proportions | *B* | **12–15%** |
| **Unit 7:** Inference for Quantitative Data: Means | *A* | **10–18%** |
| **Unit 8:** Inference for Categorical Data: Chi-Square | | **2–5%** ✗ |
| **Unit 9:** Inference for Quantitative Data: Slopes | | **2–5%** ✗ |

intersection 交集    union 并集

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \qquad P(A \cap B) = P(A) \cdot P(B|A)$$

conditional probability ('given that') : 两条件交集

given that 后条件

independent

$$P(A \cap B) = P(A) \cdot P(B)$$

$$P(A|B) = P(A|\bar{B})$$

加减不影响标准差

**Conditional Probability Formula**

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

Probability of A given B

Probability of A and B

Probability of B

$$Y = a + bx$$

$$\mu_Y = a + b\mu_x$$

$$\sigma_Y = |b| \sigma_x$$

两边同时平方变成方差    $\sigma_Y^2 = |b| \sigma_x^2$

$$\Downarrow$$

$$Var(Y) = |b|^2 Var(X)$$

X and Y are independent

if $D = X - Y$

$$E(D) = \mu_D = \mu_x - \mu_Y$$

$$\sigma_D = \sqrt{\sigma_x^2 \oplus \sigma_Y^2}$$

$$T = X + Y$$

$$\mu_T = \mu_x + \mu_Y$$

$$\sigma_T^2 = \sigma_x^2 + \sigma_Y^2 \Rightarrow \sigma_T = \sqrt{\sigma_x^2 + \sigma_Y^2}$$

$$\sqrt{n\sigma_x^2 + n\sigma_Y^2}$$

$$\sqrt{n_1 p_1 q_1 + n_2 p_2 q_2}$$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

e.g. $P(A) : 0.8 \qquad P(B) : 0.7$

$0.8 + 0.7 - 0.8 \times 0.7 = 0.94$

线性关系变量求均值

$$y = 5x + 2$$
$$\downarrow$$
$$y_1 = 5x_1 + 2$$
$$y_2 = 5x_2 + 2$$
$$y_3 = 5x_3 + 2$$

$$\frac{y_1 + y_2 + y_3}{3} = \frac{5x_1 + 2 + 5x_2 + 2 + 5x_3 + 2}{3}$$

$$\bar{y} = 5\bar{x} + 2$$

$$\boxed{\bar{y} = k\bar{x} + C}$$

已知 $\sigma x^2$

$$\sigma y^2 = \frac{(y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 + (y_3 - \bar{y})^2}{n}$$

$$= \frac{[5x_1 + 2 - (5\bar{x} + 2)]^2 + [5x_2 + 2 - (5\bar{x} + 2)]^2 + [5x_3 + 2 - (5\bar{x} + 2)]^2}{n}$$

$$= \frac{5^2 \sum (x_i - \bar{x})^2}{n}$$

$$= 5^2 \sigma x^2$$

$$\sigma y = 5 \sigma x$$

$$\sigma y^2 = k^2 \sigma x^2 \rightarrow \boxed{\sigma y = k \sigma x}$$

An analogy:
Probability: Starting with an animal, and figuring out what footprints it will make.
Statistics: seeing a footprint, and guessing the animal.

## Normal distribution as approximation to binominal 正态分布近似二项分布

$$\mu_x = np \qquad \sigma_x = \sqrt{np(1-p)}$$

条件: $np \geq 10$, $nq \geq 10$

① 求出参数 ($\mu$, $\sigma$)
② Norm. Cdf (连续性较正, $n \pm 0.5$)

### 大数定理
重复实验次数越多，事件发生的概率就越接近期望值.

### 二项分布
$$\mu = np \qquad \sigma = \sqrt{np(1-p)}$$ (当 $n = 1$ 时是伯努力分布)
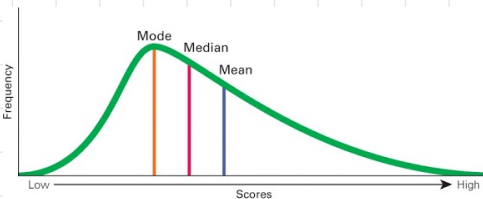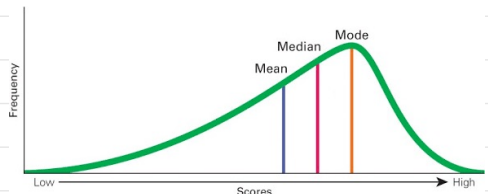
### 几何分布
$$\mu = \frac{1}{p} \qquad \sigma = \sqrt{\frac{1-p}{p^2}}$$

方差: 数据的波动大小

$$s^2 = \frac{1}{n}\left[(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2\right]$$



(a) Right-skewed distribution



(b) Left-skewed distribution

$$E(\bar{X}) = E\left(\frac{x_1 + x_2 + \cdots x_n}{n}\right)$$

$$= E\left(\frac{x_1}{n}\right) + E\left(\frac{x_2}{n}\right) + \cdots + E\left(\frac{x_n}{n}\right)$$

$$= \frac{1}{n}E(x_1) + \frac{1}{n}E(x_2) + \cdots + \frac{1}{n}E(x_n)$$

$$= \frac{1}{n} \cdot \mu \cdot n$$

$$= \mu$$

$$Var(\bar{X}) = Var\left(\frac{x_1 + x_2 + \cdots + x_n}{n}\right)$$

$$= Var\left(\frac{x_1}{n}\right) + Var\left(\frac{x_2}{n}\right) + \cdots + Var\left(\frac{x_n}{n}\right)$$

$$= \frac{1}{n^2}Var(x_1) + \frac{1}{n^2}Var(x_2) + \cdots + \frac{1}{n^2}Var(x_n)$$

$$= \frac{1}{n^2} \cdot \sigma^2 + \frac{1}{n^2} \cdot \sigma^2 + \cdots + \frac{1}{n^2} \cdot \sigma^2$$

$$= \frac{1}{n^2} \cdot \sigma^2 \cdot n$$

$$= \frac{\sigma^2}{n}$$

样本均值的方差是总体方差的 $\frac{1}{n}$

样本均值的标准差是总体标准差的 $\frac{1}{\sqrt{n}}$

样本均值的均值是总体的均值

如果本就服从正态分布,
样本数量不用很大

$$\sigma_{\bar{x}} = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}} \qquad \left[比例 \quad \sigma_{\hat{p}} = \frac{\sqrt{P(1-P)}}{\sqrt{n}}\right]$$

统计学第一定律

# Central Limit Theorem (CLT)

If $X \sim N(\mu, \sigma^2)$, then $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$.

If $X \sim$ any distribution with a mean $\mu$, and variance $\sigma^2$,

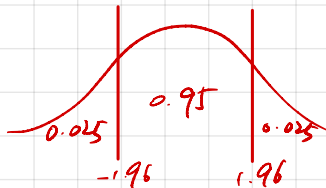then $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ given that $n$ is large.

$$\sigma_x = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$$

样本均值分布（已知总体求均值分布）

$$P\left(-1.96 < \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} < 1.96\right) = 0.95$$

（z 标注于 $\bar{x}$ 上方）

$$P\left(-1.96 \frac{\sigma}{\sqrt{n}} < \bar{x} - \mu < 1.96 \frac{\sigma}{\sqrt{n}}\right)$$

$$P\left(\mu - 1.96 \frac{\sigma}{\sqrt{n}} < \bar{x} < \mu + 1.96 \frac{\sigma}{\sqrt{n}}\right)$$



0.95

0.025    0.025

-1.96    1.96

正态分布估计标准差

$$\frac{max - min}{6} \approx \sigma$$

max: 190   min: 142

$165 - 1.96 \times 8 < \bar{x} < 165 + 1.96 \times 8$

$149.32 < \bar{x} < 180.68$

$\downarrow$

95%

正态分布标准化

$$area \rightleftharpoons z = \frac{x - \mu}{\sigma}$$ 同源

Confidence interval

marginal of error

based on the Central Limit Theorem

$$P\left(\mu - 1.96 \cdot \frac{\sigma}{\sqrt{n}} < \bar{x} < \mu + 1.96 \cdot \frac{\sigma}{\sqrt{n}}\right) = 0.95$$   已知总体均值估计样本均值

$$P\left(-1.96 \cdot \frac{\sigma}{\sqrt{n}} - \bar{x} < -\mu < 1.96 \cdot \frac{\sigma}{\sqrt{n}} - \bar{x}\right) = 0.95$$

$$P\left(1.96 \cdot \frac{\sigma}{\sqrt{n}} + \bar{x} > \mu > -1.96 \cdot \frac{\sigma}{\sqrt{n}} + \bar{x}\right) = 0.95$$

$\downarrow$

probability → 正态分布
$$P\left(\bar{x} - 1.96 \cdot \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + 1.96 \cdot \frac{\sigma}{\sqrt{n}}\right) = 0.95$$   已知样本均值估计总体的值

propotion → 比例
$$P\left(\hat{p} - z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} < P < \hat{p} + z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right) = 0.95$$   Population proportion

**t分布 σ未知**   总体标准差未知

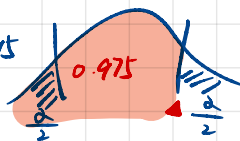$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

**t分布**           Sx是对总体 σ 的无偏估计

样本30以内   $$t = \frac{\bar{x} - \mu}{S_x / \sqrt{n}}$$

degree freedom = sample size - 1

t分布  $$P\left(\bar{x} - t_{\frac{\alpha}{2}, n-1} \cdot \frac{S}{\sqrt{n}} < \mu < \bar{x} + t_{\frac{\alpha}{2}, n-1} \cdot \frac{S}{\sqrt{n}}\right) = 0.95$$

inverse t → critical value



0.975   $\frac{\alpha}{2}$   $\frac{\alpha}{2}$

of the sample mean $\bar{x}$ is $\frac{s_x}{\sqrt{n}}$ (无系数) <span>proportion</span>

margin of error in the confidence interval for p is $ME = z \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

(有系数)

高成本的有把握 → $n \uparrow$, margin of error $\downarrow$, 置信区间变窄 / 低成本的无把握 $n \downarrow$
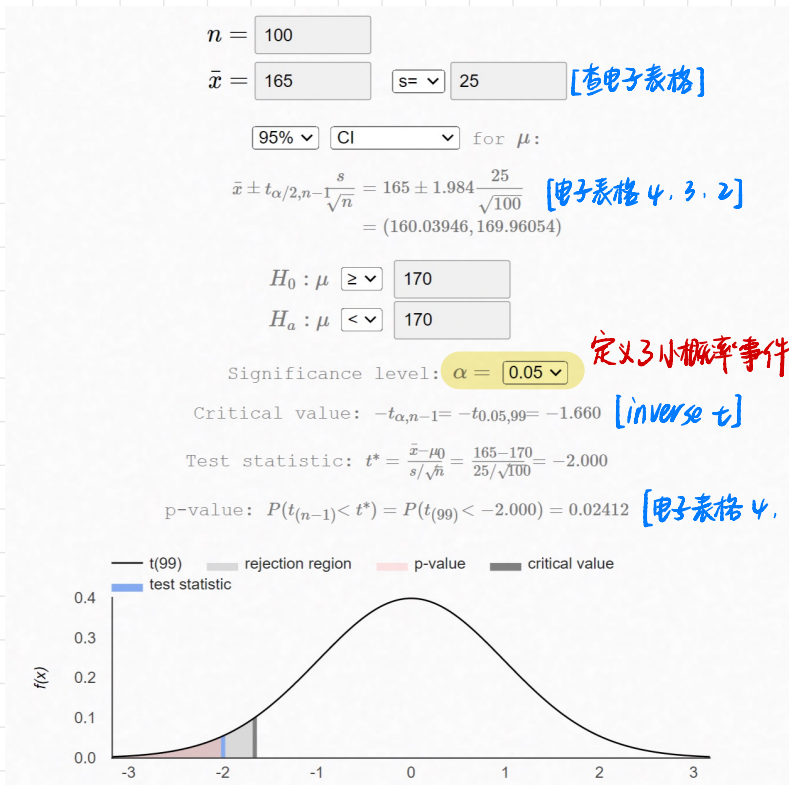
PK
① significance level 显著性水平 → 假设检验

p-value 该事件和比该事件更离谱的事件概率和 → $\begin{bmatrix} tcdf & 5,5,5 \\ 带入 t^* \end{bmatrix}$

weird

② $\begin{cases} t^* = \dfrac{\bar{x} - \mu_0}{s/\sqrt{n}} & 标准化 \\ & standard\ error \quad 样本均值标准差 \\ critical\ value & 已标准化 \quad [查表\ inverse\ t] \end{cases}$

Confidence interval

sample mean / proportion / slope $\pm$ critical value $\times$ standard error

总体均值假设检验

$$n = \boxed{100}$$

$$\bar{x} = \boxed{165} \quad \boxed{s=} \boxed{25} \quad [查电子表格]$$

$$\boxed{95\%} \quad \boxed{CI} \quad for \ \mu:$$

$$\bar{x} \pm t_{\alpha/2, n-1}\frac{s}{\sqrt{n}} = 165 \pm 1.984\frac{25}{\sqrt{100}} \quad [电子表格\ 4,3,2]$$
$$= (160.03946, 169.96054)$$

$$H_0 : \mu \boxed{\geq} \boxed{170}$$
$$H_a : \mu \boxed{<} \boxed{170}$$

Significance level: $\alpha = \boxed{0.05}$   定义3小概率事件

Critical value: $-t_{\alpha, n-1} = -t_{0.05, 99} = -1.660$   [inverse t]

Test statistic: $t^* = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{165-170}{25/\sqrt{100}} = -2.000$

p-value: $P(t_{(n-1)} < t^*) = P(t_{(99)} < -2.000) = 0.02412$   [电子表格 4,4,2]



假设检验中, Ha 决定在右尾     从单尾变双尾 P value 翻倍

双样本（样本均值差标准化）

$$\text{Test statistic} = \frac{(\bar{x}_1 - \bar{x}_2) - (\Delta_0)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$H_0$ 的值

$\Rightarrow \boxed{S_P} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$

替代 $\sigma$

样本均值差的标准差（$P_{210}$）

$\sigma$ 未知时用 $S_P$ 代替

$\begin{cases} S \begin{cases} = \\ \neq \end{cases} \\ \sigma \quad \text{样本方差} \end{cases}$

$S_P \rightarrow$ pooled    合并后的标准差

$$S_P = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}}$$

双样本差的 confidence level

$$\bar{x}_1 - \bar{x}_2 \pm t_{\frac{\alpha}{2}} \cdot \underset{\text{自由度}}{(n_1 + n_2 - 2)} \cdot S_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

双样本总体均值差假设检验



**Statistical Inference for $\mu_1 - \mu_2$**

$n_1 = 400$   $\bar{x}_1 = 171$   $\sigma_1 = 20$

$n_2 = 360$   $\bar{x}_2 = 168$   $\sigma_2 = 16$

95% CI for $\mu_1 - \mu_2$:

$$\bar{x}_1 - \bar{x}_2 \pm z_{\alpha/2}\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = 171 - 168 \pm 1.960\sqrt{\frac{20^2}{400} + \frac{16^2}{360}}$$

$$= 3.00000 \pm 2.56382$$

$$= (0.43618, 5.56382)$$

$H_0: \mu_1 - \mu_2 = 6$

$H_a: \mu_1 - \mu_2 \neq 6$

Significance level: $\alpha = 0.05$

Critical value: $\pm z_{\alpha/2} = \pm z_{0.025} = \pm 1.960$

Test statistic:

$$z^* = \frac{(\bar{x}_1 - \bar{x}_2) - (\Delta_0)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} = \frac{(171 - 168) - (6)}{\sqrt{\frac{20^2}{400} + \frac{16^2}{360}}} = -2.293$$

p-value: $2P(Z > |z^*|) = 2P(Z > 2.293) = 0.02182$

单比例假设检验 (单变量)

$$P\left(\hat{p} - z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} < p < \hat{p} + z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right) = 0.95$$

$t^*: z = \dfrac{\hat{p} - p_0 \;\;— H_0}{\sqrt{\dfrac{p_0(1-p_0)}{n}}}$

## Statistical Inference for $p$

$n =$ [ 100 ]     $\hat{p}=$ ⌄ [ 0.5 ]

Inference method: [ Wald ⌄ ]

[ 95% ⌄ ] CI for $p$:

$\hat{p} \pm z_{\alpha/2} \sqrt{\dfrac{\hat{p}(1-\hat{p})}{n}} = 0.50000 \pm 1.960 \sqrt{\dfrac{0.50000(1-0.50000)}{100}}$

$= 0.50000 \pm 0.09800$

$= (0.40200, 0.59800)$

standard error
标准误

$H_0 : p$ [ = ⌄ ] [ 0.6 ]

$H_a : p$ [ ≠ ⌄ ] [ 0.6 ]

Significance level: $\alpha =$ [ 0.05 ⌄ ]

Critical value: $\pm z_{\alpha/2} = \pm z_{0.025} = \pm 1.960$

Test statistic:

$z^* = \dfrac{\hat{p} - p_0}{\sqrt{\dfrac{p(1-p)}{n}}} = \dfrac{0.50000 - (0.6)}{\sqrt{\dfrac{0.50000(1-0.50000)}{100}}} = -2.000$

p-value: $2P(Z > |z^*|) = 2P(Z > 2.000) = 0.04550$

# 双样本比例差假设检验 (双变量)

## Statistical Inference for $p_1 - p_2$

$n_1 =$ [400]   $\hat{p}_1 =$ [0.7]

$n_2 =$ [600]   $\hat{p}_2 =$ [0.65]

[95% ▾] CI for $p_1 - p_2$:

$$\hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2}\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

$$= 0.7 - 0.65 \pm 1.960\sqrt{\frac{0.7(1-0.7)}{400} + \frac{0.65(1-0.65)}{600}}$$

$$= 0.05000 \pm 0.05893$$
$$= (-0.00893, 0.10893)$$

$H_0 : p_1 - p_2$ [= ▾] [0.12]
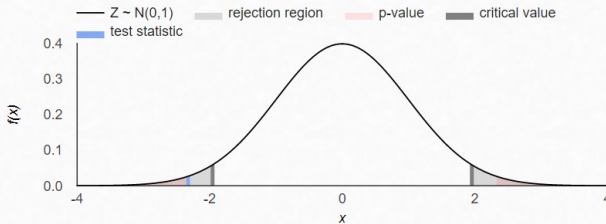
$H_a : p_1 - p_2$ [≠ ▾] [0.12]

Significance level: $\alpha =$ [0.05 ▾]

Critical value: $\pm z_{\alpha/2} = \pm z_{0.025} = \pm 1.960$

Test statistic:

$$z^* = \frac{(\hat{p}_1 - \hat{p}_2) - (\Delta_0)}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}} = \frac{(0.7 - 0.65) - (0.12)}{\sqrt{\frac{0.7(1-0.7)}{400} + \frac{0.65(1-0.65)}{600}}} = -2.328$$

p-value: $2P(Z > |z^*|) = 2P(Z > 2.328) = 0.01991$

**Paired t-test** <u>same subject</u>

test statistic $= \dfrac{\bar{d} - k}{\frac{S_d}{\sqrt{n}}}$    $\bar{d}$ : 样本差的均值

$k$ : $H_0$

$S_d$ : 样本差的标准差    (P264)

(P214) confidence interval

反证法

e.g. 设 $\sqrt{2} = \dfrac{a}{b}$ (a, b互质) → 设 $\sqrt{2}$ 为有理数

$\sqrt{2}\,b = a$

$2b^2 = a^2$

$2b^2 = (2c)^2$ → a 是偶数

$2b^2 = 4c^2$

$b^2 = 2c^2$ → b 也是偶数

∴ a, b 不可能互质

∴ $\sqrt{2}$ 不是有理数

$\begin{cases} \text{type I error : 弃真 : 放弃正确 } H_0 \\ \text{type II error : 纳伪 : 接纳错误 } H_0 \end{cases}$

**type II error**    错误的认为 $H_0$ 是对的，没有能够拒绝 $H_0$，$H_0$ 实际上是不对的

5 5 3 = critical value

area = 1 - significance level    $\mu$ : 假设分布平均值    $\sigma$ : $\dfrac{\sigma}{\sqrt{n}}$

5 5 2

upper bound : 积分积到 critical value    $\mu$ : true mean (唯一真分布的平均值)    $\sigma$ : $\dfrac{\sigma}{\sqrt{n}}$

**type I error** = significance level

不标准 (原分布)

$\dfrac{\text{critical value} - \mu_a}{\frac{\sigma}{\sqrt{n}}}$ = 标准化后 critical value (area$^{-1}$)

[5.5.3] inverse N.

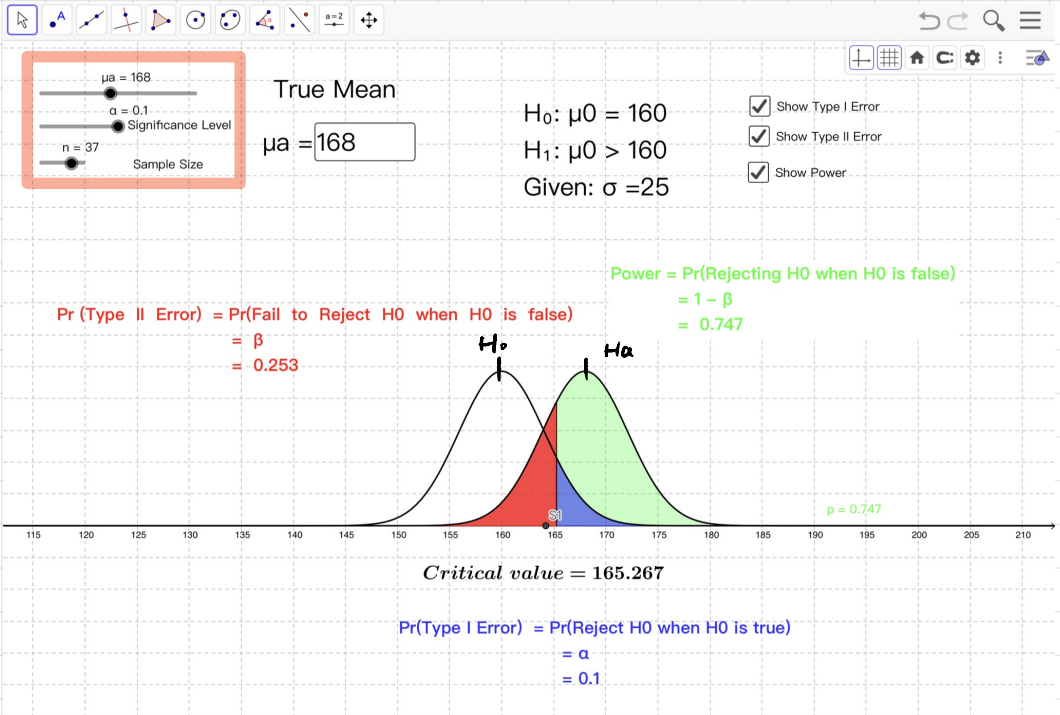| | 原假设 $H_0$ 为真 | 原假设 $H_0$ 为假 |
|---|---|---|
| 拒绝 $H_0$ | type I error ($\alpha$) | 正确决策 ($1-\beta$) |
| 不拒绝 $H_0$ | 正确决策 ($1-\alpha$) | type II error |

e.g. ① $\dfrac{166.854 - \mu a}{\frac{25}{\sqrt{36}}} = -0.524401$  $(\phi^{-1}\, 0.3)$

$\mu a = 169.039$

② $\dfrac{166.854 - \mu a}{4.16667} = -0.253347$

$\mu a = 167.91$

③ $\dfrac{166.854 - \mu a}{\frac{25}{\sqrt{36}}} = -0.582842$

$\mu a = 169.283$



$\mu_0,\ \alpha,\ n$ 对 $\beta$ ( type II error 的影响 )

两个波峰距离越小, $\beta \uparrow$     $|H_0 - Ha| \uparrow$ , power $\uparrow$

$\alpha \downarrow$ , $\beta \uparrow$

$n \downarrow$ , $\beta \uparrow$

## 最小二乘法

$\Sigma \, (y_i - y_{pi})^2$ 越小越好     $y_i$: 实际值    $y_{pi}$: 预测线上的值

平方和最小线 — 最小二乘回归线 *least-square regression line*
$$(LSRL)$$

**推**

$$\frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}{n}$$

$$= \frac{\Sigma x_i^2 + n\bar{x}^2 - 2\bar{x} \cdot n\bar{x}}{n}$$

$$= \frac{\Sigma x_i^2}{n} - (\bar{x})^2 \ \checkmark \qquad\qquad \left(\frac{\Sigma x}{n}\right)^2 = \bar{x}^2$$

$$= E(x^2) - (E(x))^2$$

$$\frac{1^2 + 2^2 + 3^2 + 4^2 + 5^2 + 6^2}{6} - (3.5)^2 = 2.916667$$

$$\sigma = \sqrt{2.916667} = 1.71 \quad [\text{表格 } 4.1.1]$$

$$\sum_{i=1}^{n} (a + bx_i - y_i)^2$$

$$= \sum_{i=1}^{n} (a^2 + b^2 x_i^2 + y_i^2 + 2abx_i - 2ay_i - 2bx_iy_i)$$

$$= \sum_{i=1}^{n} a^2 + b^2 \sum_{i=1}^{n} x_i^2 + \sum_{i=1}^{n} y_i^2 + 2ab \sum_{i=1}^{n} x_i - 2a \sum_{i=1}^{n} y_i - 2b \sum_{i=1}^{n} x_iy_i$$

$$= n \times a^2 + b^2 \sum_{i=1}^{n} x_i^2 + \sum_{i=1}^{n} y_i^2 + 2abn\bar{x} - 2an\bar{y} - 2b \sum_{i=1}^{n} x_iy_i$$

$$= n \times a^2 + 2a\left(bn\bar{x} - n\bar{y}\right) + b^2 \sum_{i=1}^{n} x_i^2 + \sum_{i=1}^{n} y_i^2 - 2b \sum_{i=1}^{n} x_iy_i$$

$$= n\left(a^2 + 2a\left(b\bar{x} - \bar{y}\right) + \left(b\bar{x} - \bar{y}\right)^2\right) - n\left(b\bar{x} - \bar{y}\right)^2 + b^2 \sum_{i=1}^{n} x_i^2 + \sum_{i=1}^{n} y_i^2 - 2b \sum_{i=1}^{n} x_iy_i$$

$$= n\left(a + b\bar{x} - \bar{y}\right)^2 - n\left(b^2\left(\bar{x}\right)^2 + \left(\bar{y}\right)^2 - 2b\bar{x}\bar{y}\right) + b^2 \sum_{i=1}^{n} x_i^2 + \sum_{i=1}^{n} y_i^2 - 2b \sum_{i=1}^{n} x_iy_i$$

$$= n\left(a + b\bar{x} - \bar{y}\right)^2 - nb^2\left(\bar{x}\right)^2 - n\left(\bar{y}\right)^2 + 2nb\bar{x}\bar{y} + b^2 \sum_{i=1}^{n} x_i^2 + \sum_{i=1}^{n} y_i^2 - 2b \sum_{i=1}^{n} x_iy_i$$

$$= \left(a + b\bar{x} - \bar{y}\right)^2 + \left(\sum_{i=1}^{n} x_i^2 - n\left(\bar{x}\right)^2\right) \times b^2 + \left(2n\bar{x}\bar{y} - 2\sum_{i=1}^{n} x_iy_i\right) \times b - n\left(\bar{y}\right)^2 + \sum_{i=1}^{n} y_i^2$$

$$\begin{cases} b = -\dfrac{\left(2n\bar{x}\bar{y} - 2\sum_{i=1}^{n} x_iy_i\right)}{2\left(\sum_{i=1}^{n} x_i^2 - n\left(\bar{x}\right)^2\right)} = -\dfrac{n\bar{x}\bar{y} - \sum_{i=1}^{n} x_iy_i}{\sum_{i=1}^{n} x_i^2 - n\left(\bar{x}\right)^2} = \dfrac{\sum_{i=1}^{n} \frac{x_iy_i}{n} - \bar{x}\bar{y}}{\sum_{i=1}^{n} \frac{x_i^2}{n} - \left(\bar{x}\right)^2} = \dfrac{\sum_{i=1}^{n} \frac{x_iy_i}{n} - \bar{x}\bar{y}}{\sum_{i=1}^{n} \frac{x_ix_i}{n} - \bar{x}\bar{x}} \\ a = \bar{y} - b\bar{x} \end{cases}$$

$$\bigstar \quad y = \frac{\sum_{i=1}^{n} \frac{x_i y_i}{n} - \bar{x}\bar{y}}{\sum_{i=1}^{n} \frac{x_i x_i}{n} - \bar{x}\bar{x}} \cdot x + (\bar{y} - b\bar{x}) \qquad \text{[表格 4, 1, 3]}$$

b

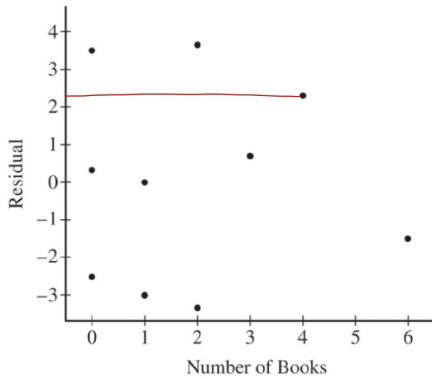$$\bar{y} = b\bar{x} + (\bar{y} - b\bar{x})$$

回归线一定经过 $(\bar{x}, \bar{y})$

$(1, 1) (2, 3) (4, 10)$

$$\frac{\frac{1\times1 + 2\times3 + 4\times10}{3} - \left(\frac{1+2+4}{3}\right) \cdot \left(\frac{1+3+10}{3}\right)}{\frac{1\times1 + 2\times2 + 4\times4}{3} - \left(\frac{1+2+4}{3}\right) \cdot \left(\frac{1+2+4}{3}\right)} \quad x + \left(\frac{1+3+10}{3} - 3.07 \cdot \frac{1+2+4}{3}\right)$$

y

b

$= 3.07 x - 2.5$

---

32. The weight, in pounds, of a full backpack and the corresponding number of books in the backpack were recorded for each of 10 college students. The resulting data were used to create the residual plot and the regression output shown below.



$y = a + bx = 10.53 + 0.53x = 10.53 + 0.53 \cdot 4 = 12.65 \qquad 12.65 + 2.2 = 14.65$

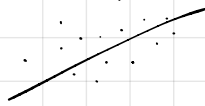| | Parameter | Estimate | Std. Err. | Alternative | DF | T-Stat | P-Value |
|---|---|---|---|---|---|---|---|
| 截距 a | Intercept | 10.53 | 1.23 | ≠ 0 | 8 | 8.57 | < 0.0001 |
| slope b | Slope | 0.53 | 0.46 | ≠ 0 | 8 | 1.15 | 0.2825 |

Which of the following values is closest to the actual weight, in pounds, of the backpack for the student who had 4 books in the backpack?
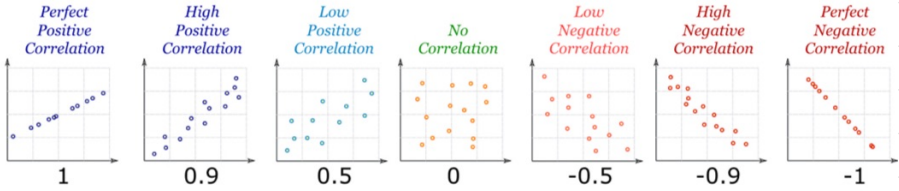
(A) 8

(B) 10

(C) 13

(D) 15

(E) 17

**r 相关系数**

$$-1 < r < 1$$

点和线的贴近程度

A correlation is assumed to be **linear** (following a line).



| Perfect Positive Correlation | High Positive Correlation | Low Positive Correlation | No Correlation | Low Negative Correlation | High Negative Correlation | Perfect Negative Correlation |
|---|---|---|---|---|---|---|
| 1 | 0.9 | 0.5 | 0 | -0.5 | -0.9 | -1 |

$$r = \frac{1}{n-1} \sum \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{x}}{s_y} \right) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$ 无单位

$$y = a + \underset{\downarrow}{b}x$$
slope

二乘回归线的斜率 / 陡峭程度

SST (total variability) — SSE (unexplained) = SSR (explained)
残余误差

$$SSE + SSR = SST$$

$$r^2 = \frac{SSR}{SST}$$     proportion of the variation that can be explained

$$SSE = 0 \Rightarrow SST = SSR \Rightarrow 相关性最强$$

$$\frac{\sqrt{\sum (y_i - \bar{y})^2}}{\sqrt{\sum (x_i - \bar{x})^2}} = \frac{s_y}{s_x}$$     计算式

$$b = r \cdot \frac{s_y}{s_x}$$   b, r  关系式

# 最小二乘回归线的置信区间和假设检验



Regression

$\hat{y}_2 = a_2 + b_2 x$

$\hat{y} = a + bx$

Height / Shoe size

$\hat{y} = \alpha + \beta x$

critical value

$b_2 \pm t^* \cdot SE_b$ ← Confidence interval

★ $H_0: \beta = 0$
$H_a: \beta \neq 0$ ← Hypothesis test

Regression

$\hat{y}_2 = a_2 + b_2 x$

$\hat{y} = a + bx$

Height / Shoe size

$\hat{y} = \alpha + \beta x$

Conditions for Inference

Linear — Actual linear relationship between x & y

① Independence
✓ Normal
✓ Equal variance
✓ Random

---

Musa is interested in the relationship between hours spent studying and caffeine consumption among students at his school. He randomly selects <u>20 students</u> at his school and records their caffeine intake (mg) and the amount of time spent studying in a given week. Here is computer output from a least-squares regression analysis on his sample:

| Predictor | Coef | SE Coef | T | P |
|---|---|---|---|---|
| Constant $a =$ | 2.544 | 0.134 | 18.955 | 0.000 |
| Caffeine $\beta =$ | 0.164 | 0.057 | 2.862 | 0.010 |
| S = 1.532 | R-sq = 60.0% | SE | $t^*$ | p-value |

Time Study slope = 0.164

Caffeine

Assume that all conditions for inference have been met.

**What is the 95% confidence interval for the slope of the least squares regression line?**

$0.164 \pm t^* \cdot 0.057$

$df = 20 - 2 = \boxed{18}$

---

Jian obtained a random sample of data on how long it took each of 24 students to complete a timed reaction game and a timed memory game. He noticed a positive linear relationship between the times on each task. Here is computer output on the sample data:

Population — samples 24 data points

$\hat{y} = \alpha + \beta x$

$H_0: \beta = 0$
$H_a: \beta > 0$

$\hat{y} = \alpha + \beta x$

### Summary statistics

| Variable | n | Mean | StDev | SE Mean |
|---|---|---|---|---|
| $x$ = reaction time | 24 | 0.398 | 0.133 | 0.027 |
| $y$ = memory time | 24 | 43.042 | 8.554 | 1.746 |

### Regression: memory vs. reaction

| Predictor | Coef | SE Coef |
|---|---|---|
| Constant $a =$ | 37.200 | 5.579 |
| Reaction $b =$ | 14.686 | 13.329 |
| S = 8.515 | R-sq = 5.2% | |

Assume that all conditions for inference have been met.

Calculate the test statistic that should be used for testing a null hypothesis that the population slope is actually 0?

$z = \dfrac{b - \beta_0}{\sigma_b}$

$t = \dfrac{b - \beta_0}{SE_b} = \dfrac{14.686 - 0}{13.329}$

---

Alicia took a random sample of mobile phones and found a positive linear relationship between their processor speeds and their prices. Here is computer output from a least-squares regression analysis on her sample:

### Regression: Price vs. speed

| Predictor | Coef | SE Coef | T | P |
|---|---|---|---|---|
| Constant | 127.092 | 57.507 | 2.210 | 0.032 |
| Speed | 6.084 | 2.029 | 2.999 | 0.004 |

Price / Proc Speed $\hat{y} = \alpha + \beta x$

Price / Proc. Speed $\hat{y} = \alpha + \beta x$

p-value = $P(t \geq 2.999) = 0.002$

Alicia wants to test $H_0: \beta = 0$ vs. $H_a: \beta > 0$. Assume that all conditions for inference have been met.

**At the $\alpha = 0.01$ level of significance, is there sufficient evidence to conclude a positive linear relationship between these variables for all mobile phones? Why?**

$0.002 < 0.010 \Rightarrow$ reject $H_0$

Yes, because P-value < $\alpha$ ⇒ reject $H_0$ ⇒ suggests

---

Hashem obtained a random sample of students and noticed a positive linear relationship between their ages and their backpack weights. A <u>95% confidence interval</u> for the slope of the regression line was $0.39 \pm 0.23$.

Hashem wants to use this interval to test $H_0: \beta = 0$ vs. $H_a: \beta \neq 0$ at the $\alpha = 0.05$ level of significance. Assume that all conditions for inference have been met.

95% confidence interval: $[0.16, 0.62]$

Assuming $H_0$ true, we are in $\leq 5\%$ situations where $\beta$ not overlap with 95% interval

↓

reject $H_0 \Rightarrow$ suggest $H_a: \beta \neq 0$

↓

there is a non-zero linear relationship between ages & backpack weights.

**卡方分布** (inference for categorical data)

**独立性检验 & 同质性检验**

association $\begin{cases} \text{independence} \\ \text{homogeneity} \rightarrow \text{proportion} \end{cases}$

$H_0$: independent
$H_a$: not independent

$df = (r-1)(c-1)$    $TS = \chi^2 = \Sigma \frac{(O_i - E_i)^2}{E_i}$    expected value = $\frac{R_i \cdot C_i}{n}$

① 萃 7, 1, 1    control var.    enter 输入矩阵
                                    ↓
                                  store

② 萃 6, 7, 8

expected value : var. expect.
         (4)

**拟合优度检验**

有一标准分布，检验抽样是否符合

$df = k - 1$

expected value
电子表格   再加一列, =, var. pro, × 总数, enter

电子表格   4, 4, 7    compare $x^2$ with critical value
                         5, 5, 9    卡方检验都是单尾

in all conditions, $E_i \geq 5$, then $\Sigma \left( \frac{(O_i - E_i)^2}{E_i} \right) \sim x^2$ is a good approximation.

sample size   一般小于等于总体的 10% → $n \leq \frac{1}{10} N$

$H_0$: 两者一致 (TS 小, P value 大)
$H_a$: 两者不一致 (TS 大, P value 小)

Ha 对应小概率事件 ⇄ 大 $x^2$ 值 ⇄ 变量比例高度不一致

$$\begin{bmatrix} 60 & 20 \\ 60 & 30 \\ 60 & 90 \end{bmatrix}$$

$x^2 = 31.5$

$P\ value = 1.45 \times 10^{-7}$

拒绝 $H_0$，接受 Ha    一定不独立

$$\begin{bmatrix} 60 & 20 \\ 60 & 20 \\ 60 & 20 \end{bmatrix}$$

$x^2 = 0$

$P\ value = 1$

无法拒绝 $H_0$    可能独立

~~希望~~小概率事件发生，以此拒绝 $H_0$，接受 Ha.

在 $H_0$ 的条件下，

小概率事件一旦发生，说明 $H_0$ 错误 → 所以 Ha 正确
($P\ value < 0.05$)

小概率事件不发生，$H_0$ 有可能是对的 → 无法拒绝 $H_0$
($P\ value > 0.05$)

在卡方检验中，$H_0$ 表示相类似 / 无差异 / 同比例 / 相互独立.
                        no association

所有 $P\ value$ 都是基于 $H_0$ 成立

**P value**：$H_0$ 正确的条件下，其他更极端事件发生概率

| | 数学好 | 数学不好 | | | |
|---|---|---|---|---|---|
| 左 | 70 | 30 | | 70 | 30 |
| 右 | 30 | 70 | | 70 | 30 |

$x^2 = 32$     $\left( \dfrac{(70-50)^2}{50} + \dfrac{(30-50)^2}{50} \right) \times 2$     $x^2 = 0$

$(8 + 8) \times 2 = 32$

# 统计问题定性分析

1. confidence interval    vs.  hypothesis test  ( prediction / Chi - square )
    $a \leq x \leq b$                statement  with  $\leq$ , $\geq$ , $=$

2. proportion  vs.  mean
   · 男生人数占比          · 身高
   (女生人数占比)          · 体重
   · 得新冠的患病率      · 血压

3. independent  vs.  dependent
   (unpaired  vs.  paired )

4. one  sample  vs.  two  samples

5. categorical  vs.  quantitative

In 95% of all sample condition, the method  would  yeild  an
interval  captures  the  true  parameter  value.

1. 笔记
2. flashcards  quizelet
3. 统汇问题定性分析
   http://www.ltcconline.net/greenl/java/Statistics/catStatProb/categorizingStatProblemsJavaScript.html

4. 定量计算 ( TI - Nspire )
5. 官方单元真题
6. 往年真题
7. 必考   ① 双变量
          ② 概率
      ★ ③  confidence  interval
      ★ ④ 假设检验
          ⑤ 卡方 / slope